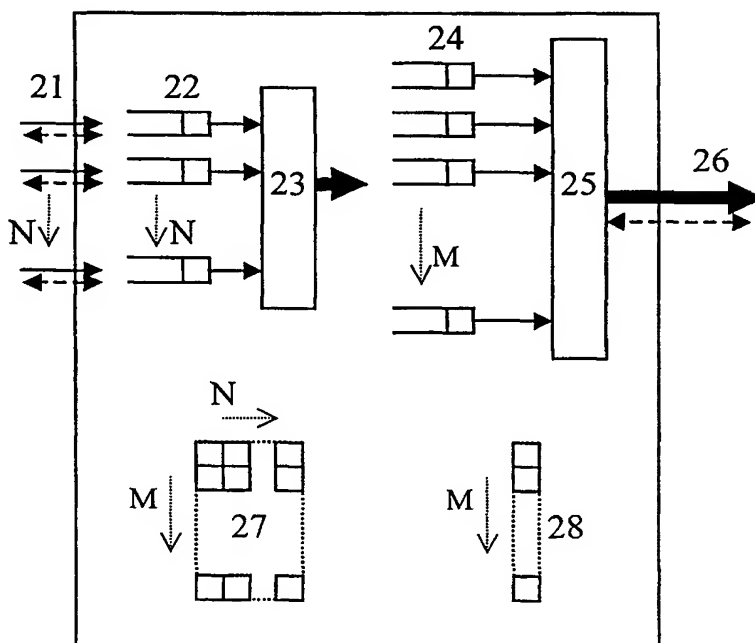




INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

| | | | | |
|---|---|---|---|---|
| (51) International Patent Classification ⁷ : H04L 12/56 | A1 | (11) International Publication Number: WO 00/38376 (43) International Publication Date: 29 June 2000 (29.06.00) | | |
| <table style="width: 100%; border: none;"> <tr> <td style="width: 50%; vertical-align: top; padding: 5px;"> (21) International Application Number: PCT/GB99/04007 (22) International Filing Date: 1 December 1999 (01.12.99) (30) Priority Data: 9828143.9 22 December 1998 (22.12.98) GB (71) Applicant (for all designated States except US): POWER X LIMITED [GB/GB]; Stafford Court, 145 Washway Road, Sale, Manchester M33 7PE (GB). (72) Inventors; and (75) Inventors/Applicants (for US only): PIEKARSKI, Marek, Stephen [GB/GB]; 20 Gawsworth Road, Macclesfield, Cheshire SK11 8UE (GB). JOHNSON, Ian, David [GB/GB]; 11 Seel Street, Mosley OL5 0EW (GB). (74) Agents: McNEIGHT, David, Leslie et al.; McNeight & Lawrence, Regent House, Heaton Lane, Stockport, Cheshire SK4 1BS (GB). </td> <td style="width: 50%; vertical-align: top; padding: 5px;"> (81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG). Published <i>With international search report.</i> </td> </tr> </table> | | | (21) International Application Number: PCT/GB99/04007 (22) International Filing Date: 1 December 1999 (01.12.99) (30) Priority Data: 9828143.9 22 December 1998 (22.12.98) GB (71) Applicant (for all designated States except US): POWER X LIMITED [GB/GB]; Stafford Court, 145 Washway Road, Sale, Manchester M33 7PE (GB). (72) Inventors; and (75) Inventors/Applicants (for US only): PIEKARSKI, Marek, Stephen [GB/GB]; 20 Gawsworth Road, Macclesfield, Cheshire SK11 8UE (GB). JOHNSON, Ian, David [GB/GB]; 11 Seel Street, Mosley OL5 0EW (GB). (74) Agents: McNEIGHT, David, Leslie et al.; McNeight & Lawrence, Regent House, Heaton Lane, Stockport, Cheshire SK4 1BS (GB). | (81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG). Published <i>With international search report.</i> |
| (21) International Application Number: PCT/GB99/04007 (22) International Filing Date: 1 December 1999 (01.12.99) (30) Priority Data: 9828143.9 22 December 1998 (22.12.98) GB (71) Applicant (for all designated States except US): POWER X LIMITED [GB/GB]; Stafford Court, 145 Washway Road, Sale, Manchester M33 7PE (GB). (72) Inventors; and (75) Inventors/Applicants (for US only): PIEKARSKI, Marek, Stephen [GB/GB]; 20 Gawsworth Road, Macclesfield, Cheshire SK11 8UE (GB). JOHNSON, Ian, David [GB/GB]; 11 Seel Street, Mosley OL5 0EW (GB). (74) Agents: McNEIGHT, David, Leslie et al.; McNeight & Lawrence, Regent House, Heaton Lane, Stockport, Cheshire SK4 1BS (GB). | (81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG). Published <i>With international search report.</i> | | | |
| (54) Title: DISTRIBUTED HIERARCHICAL SCHEDULING AND ARBITRATION FOR BANDWIDTH ALLOCATION | | | | |
| (57) Abstract <p>A scheduling and arbitration arrangement is provided for use in a digital switching system comprising a central switch under the direction of a master control to provide the cross-connections between a number of high-bandwidth input and output ports. On the ingress side of the switch are provided a number of ingress multiplexers, one for each high-bandwidth input port, whilst on the egress side of the switch are a number of egress multiplexers, one for each high-bandwidth output ports. Each ingress multiplexer includes a set of N input queues serving N low-bandwidth data sources and a set of M virtual output queues, one for each low-bandwidth output data source. The scheduling and arbitration arrangement includes three bandwidth allocation tables. One of these, the ingress port table, is associated with the input queues having NxM entries, each arranged to define the bandwidth allocation for a particular virtual output queue. A second table is the egress port table, associated with the virtual output queues having M entries, each arranged to define the bandwidth allocation of a high-bandwidth port of the central switch to a virtual output queue. The third table is the central allocation table, located in the master control and having $(M/N)^2$ entries, each of which specifies the weights allocated to each possible connection through the central switch.</p> | | | | |



FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

| | | | | | | | |
|----|--------------------------|----|--|----|--|----|--------------------------|
| AL | Albania | ES | Spain | LS | Lesotho | SI | Slovenia |
| AM | Armenia | FI | Finland | LT | Lithuania | SK | Slovakia |
| AT | Austria | FR | France | LU | Luxembourg | SN | Senegal |
| AU | Australia | GA | Gabon | LV | Latvia | SZ | Swaziland |
| AZ | Azerbaijan | GB | United Kingdom | MC | Monaco | TD | Chad |
| BA | Bosnia and Herzegovina | GE | Georgia | MD | Republic of Moldova | TG | Togo |
| BB | Barbados | GH | Ghana | MG | Madagascar | TJ | Tajikistan |
| BE | Belgium | GN | Guinea | MK | The former Yugoslav Republic of Macedonia | TM | Turkmenistan |
| BF | Burkina Faso | GR | Greece | | | TR | Turkey |
| BG | Bulgaria | HU | Hungary | ML | Mali | TT | Trinidad and Tobago |
| BJ | Benin | IE | Ireland | MN | Mongolia | UA | Ukraine |
| BR | Brazil | IL | Israel | MR | Mauritania | UG | Uganda |
| BY | Belarus | IS | Iceland | MW | Malawi | US | United States of America |
| CA | Canada | IT | Italy | MX | Mexico | UZ | Uzbekistan |
| CF | Central African Republic | JP | Japan | NE | Niger | VN | Viet Nam |
| CG | Congo | KE | Kenya | NL | Netherlands | YU | Yugoslavia |
| CH | Switzerland | KG | Kyrgyzstan | NO | Norway | ZW | Zimbabwe |
| CI | Côte d'Ivoire | KP | Democratic People's Republic of Korea | NZ | New Zealand | | |
| CM | Cameroon | | | PL | Poland | | |
| CN | China | KR | Republic of Korea | PT | Portugal | | |
| CU | Cuba | KZ | Kazakstan | RO | Romania | | |
| CZ | Czech Republic | LC | Saint Lucia | RU | Russian Federation | | |
| DE | Germany | LI | Liechtenstein | SD | Sudan | | |
| DK | Denmark | LK | Sri Lanka | SE | Sweden | | |
| EE | Estonia | LR | Liberia | SG | Singapore | | |

Distributed Hierarchical Scheduling and Arbitration
for Bandwidth Allocation

The present invention relates to data switching systems and is more particularly concerned with the scheduling and arbitration arrangements for such systems.

The continual growth of demand for manageable bandwidth in networks requires the development of new techniques in switch design which decouples the complexity of control from the scale of the port count and aggregate bandwidth. This invention describes a switch architecture and a set of methods which provide the means by which switches of arbitrary size may be constructed whilst maintaining the ability to allocate guaranteed bandwidth to each possible connection through the switch. A digital switch is used to route data streams from a set of source components to a set of destination components. A cell-based switch operates on data which is packetised into streams of equal size cells. In a large switch the routing functions may be implemented hierarchically, that is sets of lower bandwidth ports are aggregated into a smaller number of higher bandwidth ports which are then interconnected in a central switch.

It is an object of the present invention to provide a bandwidth allocation arrangement which may be used in such a hierarchical switch.

According to the present invention there is provided a scheduling and arbitration process for use in a digital data switching arrangement of the type in which a central switch under the direction of a master control provides the cross-connections between a number of high-bandwidth ports to which are connected on the ingress side of the central switch a number of ingress multiplexers, one for each high-bandwidth input port and on the egress side a number of egress multiplexers, one for each high-bandwidth output port, each ingress multiplexer including a set of N input queues serving N low-bandwidth data sources and a set of M virtual output queues serving M low-bandwidth output data sources, characterised in that the scheduling and arbitration arrangement includes three bandwidth allocation tables, an ingress port table associated with the input queues and having NxM entries each arranged to define the bandwidth for a particular virtual output

queue, an egress port table associated with the virtual output queues and having M entries each arranged to define the bandwidth allocation of a high-bandwidth port of the central switch to a virtual output queue and a central allocation table located in the master control and having $(M \times N)^2$ entries each of which specifies the weights allocated to each possible connection through the central switch.

According to a feature of the invention there is provided a scheduling and arbitration process in which the scheduling of the input queues is performed in accordance with an N-way weighted round robin.

According to a further feature of the invention there is provided an implementation of the N-way weighted round robin by an $N \cdot (2^w - 1)$ -way unweighted round robin where w is the number of bits defining a weight using a list constructed by interleaving N words of $(2^w - 1)$ bits each, with w_n 1's in a word, where w_n is the weight of the queue n .

The invention, together with its various features, will be more readily understood from the following description of one embodiment, which should be read in conjunction with the accompanying drawings. In the drawings:-

Figure 1 shows a simplified form of a data switch,

Figure 2 shows an egress multiplexer,

Figure 3 shows the weighted round robin arbiter for use in the egress multiplexer,

Figure 4 shows the partitioning of the round robin arbiter,

Figure 5 shows the operation of the round robin arbiter,

Figure 6 shows the allocations for a 4-port interconnect with 3 bit weights and

Figure 7 shows a block diagram of a small switch based on the principles of the invention.

Referring now to Figure 1, this shows a schematic diagram of a hierarchical switch. The central interconnect 1 provides the cross-connections between a number of high-bandwidth ports. A set of multiplexers 2 on the ingress side and demultiplexers 3 on the egress side provides the aggregation function between the low and high-bandwidth

ports. The low bandwidth ports provide connections from the switch to the data sources 4 on the ingress side and the data destinations 5 on the egress side. In practice, a switch is required to support full-duplex ports, so that an ingress multiplexer and its corresponding demultiplexer may be considered a single full-duplex device which will be hereafter termed a “router” Typically the data switch may be of the type disclosed in out co-pending patent application No. PCT/GB99/03748.

It should be noted that the central interconnect 1 may itself be a hierarchical switch, that is the methods described may be applied to switches with an arbitrary number of hierarchical levels.

The aim of these methods is to provide a mechanism whereby the data stream from the switch to a particular destination, which comprises a sequence of cells interleaved from various data sources, may be controlled such that predetermined proportions of its bandwidth are guaranteed to cells from each data source.

Figure 2 shows the architecture of an ingress multiplexer. An ingress multiplexer receives a set of data streams from the data sources via a set of low-bandwidth input ports. Each data stream is a sequence of equal size cells (that is, an equal number of bits of data). A set of N low-bandwidth ports 21 each fills one of the N input queues 22. An ingress control unit 23 extracts the destination address from each of the cells in the input queues and transfers them into a set of M virtual output queues 24. There is one virtual output queue for each low-bandwidth output port in the switch. The ingress multiplexer also contains an interconnect link control unit 25 which implements this function by scheduling cells from the virtual output queues 24 across the high-bandwidth link 26 to the central interconnect 1 according to an M-entry egress table 28.

In addition to the data flow indicated by the arrows in Figure 1, there is also a flow of backpressure or flow-control information associated with each of the data flows. This control flow is indicated in Figure 2 by dashed arrows. The ingress multiplexer contains an NxM-entry ingress port table 27, which defines how its bandwidth to a particular egress port (via a particular virtual output queue) is distributed across the input ports. This table is used by the ingress control unit 23 to determine when (and to what degree) to

exert backpressure to the data source resolved down to an individual virtual output queue.

The ingress multiplexer 2 of Figure 1 sends control information to the central interconnect 1 indicating the state of the virtual output queues in the form of "connection requests". The central interconnect responds with a sequence of connections which it will establish between the ingress and egress routers. These are "connection grants". The ingress multiplexer 2 must now allocate the bandwidth to each egress demultiplexer 3 provided by the central interconnect 1 across the virtual output queues associated with each egress demultiplexer.

The deterministic scheduling function of the interconnect link control unit 25 may be defined as a weighted round robin (WRR) arbiter. The interconnect link control unit 25 receives a connection grant to a particular egress demultiplexer 3 from the central interconnect 1 and must select one of the N virtual output queues associated with that egress demultiplexer. This may be implemented by expanding the N-way WRR shown in Figure 3a) into an $(N \cdot (2^W - 1))$ - way unweighted round robin as shown in Figure 3b), where W equals the number of bits necessary to define the weight, such that if a queue has a weight of w , then it is represented as $(w-1)$ entries in the unweighted round robin list. For example, with 4-bit weights, a 4-way weighted round robin expands to a 60-way unweighted round robin.

In order to optimise the service intervals to the queues under all weighting conditions, the entries in the unweighted round robin list are distributed such that for each weight the entries are an equal number of steps apart plus or minus one step. Table 1 below shows an example of such an arrangement of 3-bit weights:

| \underline{w}_n | \underline{e}_n |
|-------------------|-------------------|
| 1 | 1000000 |
| 2 | 1000100 |
| 3 | 1001010 |
| 4 | 1010101 |
| 5 | 1011011 |
| 6 | 1110111 |
| 7 | 1111111 |

In the system described, the arbiter must select one of the nine queues with 4-bit weights, that is 8 virtual output queues as described above and a multicast queue. This expands to a 135-entry unweighted round robin. The implementation of a large unweighted round robin arbiter may be achieved without resorting to a slow iterative shift-and-test method by the technique of “divide and conquer”, that is the 135-entry round robin is segmented into 9 sections of 16-entry round robins, each of which may be implemented efficiently with combinational logic (9 x 16 provides up to 144 entries, so that the multicast queue of up to 24 entries may actually be allocated more bandwidth than an individual unicast queue of up to 15 entries).

Figure 4 illustrates the partitioning of the round robin arbiter. The sorter 41 separates the request vector V (144 bits) into 9 sections of 16-bit vectors, v0 to v8. It also creates nine pointers p0 to p8 for each of the 16-bit round robin blocks 42. The block which corresponds to the existing pointer (which has been saved in register 44) is given a “1” at the corresponding bit location, whilst the other blocks are given dummy pointers initialised to location zero. Each 16-bit round robin block now finds the next “1” in its input vector and outputs its location (g) whether it has to wrap round (w) and whether it has found a “1” in its vector (f). A selector 43 is now able to identify the block which has found the “1” corresponding to the next “1” in the original 135-bit vector given a signal (s) from the sorter 41. This specifies which round robin block had the original pointer position. The selector 43 is itself a round robin function which may be implemented as a combinational logic function

“find the next block starting at s which has w=false and f=true
(if not found, select s)”.

Figure 5 shows an example of the above process, but for a smaller configuration for clarity. In the example, V = 12 bits, p = 4 bits, v0 - 2 = 2 bits and g0 - 2 = 2 bits. Figure 5 depicts the process performed by Figure 4 and at 51 defines the expanded current pointer (P) and the request vector (V) at 52. The sorter 41 produces segmented vectors (v) and segmented pointers (p) where the blocks marked * are dummies. The segmented results (g) of the round robin are shown at 55 whereas the results of the

selector process 43 is shown at 56, defining the expanded next pointer (P).

The central interconnect 1 provides the cross-connect function in the switch. The bandwidth allocation in the central interconnect is defined by an $(M/N)^2$ -entry central allocation table, which specifies the weights allocated to each possible connection through the central interconnect (the central interconnect has M/N high-bandwidth ports). The central allocation table contains P^2 entries, where $P=(M/N)$. Each entry w_{ie} defines the weights allocated to the connection from high-bandwidth port i to high-bandwidth port e . However, not all combinations of entries constitute a self-consistent set, that is the allocations as seen from the outputs could contradict the allocations as seen from the inputs. A set of allocations is only self-consistent if the sums of weights at each output and input are equal. Figure 6 shows a self-consistent set (a) and a non-self-consistent set (b) of allocation for a 4-port interconnect with 3-bit weights. Inputs are shown at IP and outputs at OP, with the sum designated as Σ . Assuming that the central allocation table has a self-consistent set of entries, it is possible to define the bandwidth allocation to a link between input port i and output port e with weight w_{ie} as p_{ie} , where:

$$p_{ie} = \frac{w_{ie}}{\sum_{n=0}^{(p-1)} w_{in}}$$

The egress port table defines how the bandwidth of a high-bandwidth port to the central interconnect 1 is allocated across the virtual output queues. There is no issue with self-consistence as all possible entries are self-consistent, so that the bandwidth allocation for a virtual output queue v with weight w_v is given by:

$$p_f = \frac{w_v}{\sum_{n=0}^{(N-1)} w_n}$$

Similarly, the ingress port table entries give the bandwidth allocation of a virtual output

queue to the ingress ports with weight w_f is given by:

$$p_f = \frac{w_f}{\sum_{n=0}^{(N-1)} w_n}$$

Therefore the proportion of bandwidth at an egress port v allocated to an ingress port f is given by:

$$p_{fv} = p_f \cdot p_v \cdot p_{ie}$$

In a switch which is required to maintain strict bandwidth allocation between ports (such as an ATM switch), the tables are set up via a switch management interface from a connection admission and control processor. When the connection admission and control processor has checked that it has the resources available in the switch to satisfy the connection request, then it can modify the ingress port table, the egress port table and the central allocation table to reflect the new distribution of traffic through the switch.

In contrast, a switch may be required to provide a “best effort” service. In this case the table entries are derived from a number of local parameters. Two such parameters are the length l_v of the virtual output queue v and the urgency u_v of the virtual output queue. urgency is a parameter which is derived from the headers of the cells entering the queue from the ingress ports.

A switch may be implemented which can satisfy a range of requirements (including the two above) by defining a weighting function which “mixes” a number of scheduling parameters to generate the table entries in real time according to a set of “sensitivities” to length, urgency and pseudo-static bandwidth allocation. (s_l, s_w, s_s). The requirement on the function are that it should be fast and efficient, since multiple instances occur in the critical path of a switch. In the system described the weighting function has the form:

$$w_v = \left\{ \frac{l_v^2}{2^{(1/sl)}} + \frac{p_v}{2^{(1/ss)}} + \frac{u_v}{2^{(1/su)}} \right\} \cdot (1-b_v)$$

where b_v is the backpressure applied from the egress multiplexer,

w_v is the weight of the queue as applied to the scheduler, and

p_v is a pseudo-static bandwidth allocation, such as an egress port table.

Despite the apparent complexity of this function, it may be implemented exclusively with an adder, multiplexers and small lookup tables, thus meeting the requirement for speed and efficiency. Features of this weighting function are that, for $s_l = 1.0$, $s_s = 0.0$ and $s_u = 0.0$, bandwidth is allocated locally purely on the basis of queue length, with a non-linear function, so that the switch always attempts to avoid queues overflowing. When $s_l = 0.0$, $s_s = 1.0$ and $s_u = 0.0$, bandwidth is allocated purely on the basis of pseudo-static allocations as described above. Finally, when $s_l = 0.0$, $s_s = 1.0$ and $s_u = 0.5$, bandwidth is allocated on the basis of pseudo-static allocation but a data source is allowed to “push” some data harder, when the demand arises, by setting the urgency bit in the appropriate cell headers.

Figure 7 is a block diagram of a small switch based on the above principles, showing the correct number of queues, tables and table entries. In Figure 7 there are two ingress routers 71 and 72, a central cross-bar switch 73 and two egress routers 74 and 75. Each ingress router has two low-bandwidth input ports, A and B for router 71 and ports C and D for router 72. As mentioned previously, each ingress router has an ingress port table such as 77 for router 72 and an egress port table such as 78, whereby the central switch 73 has a central allocation table 79. Assuming that each low-bandwidth port may transport 1Gbps of traffic, each high-bandwidth link may carry 2 Gbps and the switch is required to guarantee the following bandwidth allocations:

| Flow bandwidth (Gbps) | Destination Port | | | |
|--------------------------|------------------|-----|-----|-----|
| | A | B | C | D |
| A | 0.5 | 0.1 | 0.1 | 0.2 |
| B | 0.2 | 0.2 | 0.2 | 0.2 |
| C | - | 0.5 | - | 0.2 |
| D | 0.1 | 0.1 | 0.6 | 0.2 |

then the ingress port table such as 77, egress port table such as 78 and central allocation table 79 would be set up by the connection admission and control processor with the following 4- bit values (note here that there will be rounding errors due to the limited resolution of the 4-bit weights):

Ingress Port Table
(in router 71)

| | Source | |
|---|--------|---|
| | A | B |
| A | 15 | 6 |
| B | 3 | 6 |
| C | 3 | 6 |
| D | 6 | 6 |

Ingress Port Table
(in router 72)

| | Source | |
|---|--------|----|
| | A | B |
| A | 0 | 3 |
| B | 15 | 3 |
| C | 0 | 15 |
| D | 6 | 5 |

Egress Port Table
(in router 71)

| | Source |
|---|--------|
| | AB |
| A | 15 |
| B | 6 |
| C | 6 |
| D | 8 |

Egress Port Table
(in router 72)

| | Source |
|---|--------|
| | CD |
| A | 2 |
| B | 12 |
| C | 12 |
| D | 12 |

Central Allocation Table

| Source | Destination Router | |
|--------|--------------------|----|
| | AB | CB |
| | CD | |
| AB | 15 | 10 |
| CD | 10 | 15 |

CLAIMS

1. A scheduling and arbitration process for use in a digital data switching arrangement of the type in which a central switch under the direction of a master control provides cross-connections between a number of high-bandwidth ports to which ingress multiplexers are connected on the ingress side of the central switch, one for each high-bandwidth input port and to which egress multiplexers are connected on the egress side of the central switch, one for each high-bandwidth output port, each ingress multiplexer including a set of N input queues serving N low-bandwidth data sources and a set of M virtual output queues, one for each low-bandwidth output data source, characterised in that the scheduling and arbitration arrangement includes three bandwidth allocation tables, an ingress port table associated with the input queues having NxM entries each arranged to define the bandwidth allocation for a particular virtual output queue, an egress port table associated with the virtual; output queues having M entries each arranged to define the bandwidth allocation of a high-bandwidth port of the central switch to a virtual output queue, and a central allocation table located in the master control and having $(M/N)^2$ entries each of which specifies the weights allocated to each possible connection through the central switch.
2. A scheduling and arbitration arrangement as claimed in Claim 1 characterised in that the scheduling of input queues is performed in accordance with a N-way weighted round robin, where N equals the number of input queues.
3. A scheduling and arbitration arrangement as claimed in Claim 2 characterised in that the N-way weighted round robin is implemented by an $N \cdot (2^w - 1)$ - way unweighted round robin, where w is the number of bits defining a weight, using a request vector list constructed by interleaving N words of $(2^w - 1)$ bits each, with w_n "1"s in a word, where w_n is the weight of the queue n.
4. A scheduling and arbitration arrangement as claimed in Claim 3 characterised in that the request vector list is separated into a plurality of 16-bit round robin blocks, one for each queue in a multiplexer, a pointer being created for each round robin block with the block corresponding to the existing pointer having a "1" at the corresponding bit

position whilst all the other pointers are initialised to zero and each round robin block is activated to identify the block which has found the next "1" in the request vector list.

5. A scheduling and arbitration arrangement as claimed in Claim 1 characterised in that the ingress port table, the egress port table and the central allocation table are all programmed from an external source.

6. A scheduling and arbitration arrangement as claimed in Claim 5 characterised in that the external source uses local parameters defining the length of the virtual output queue and the urgency of the virtual output queue.

7. A scheduling and arbitration arrangement as claimed in Claim 6 characterised in that the external source uses a set of sensitivities relating to the length, urgency and pseudo-static bandwidth allocation.

8. A digital switching arrangement characterised by a scheduling and arbitration arrangement as claimed in any one of the preceding claims.

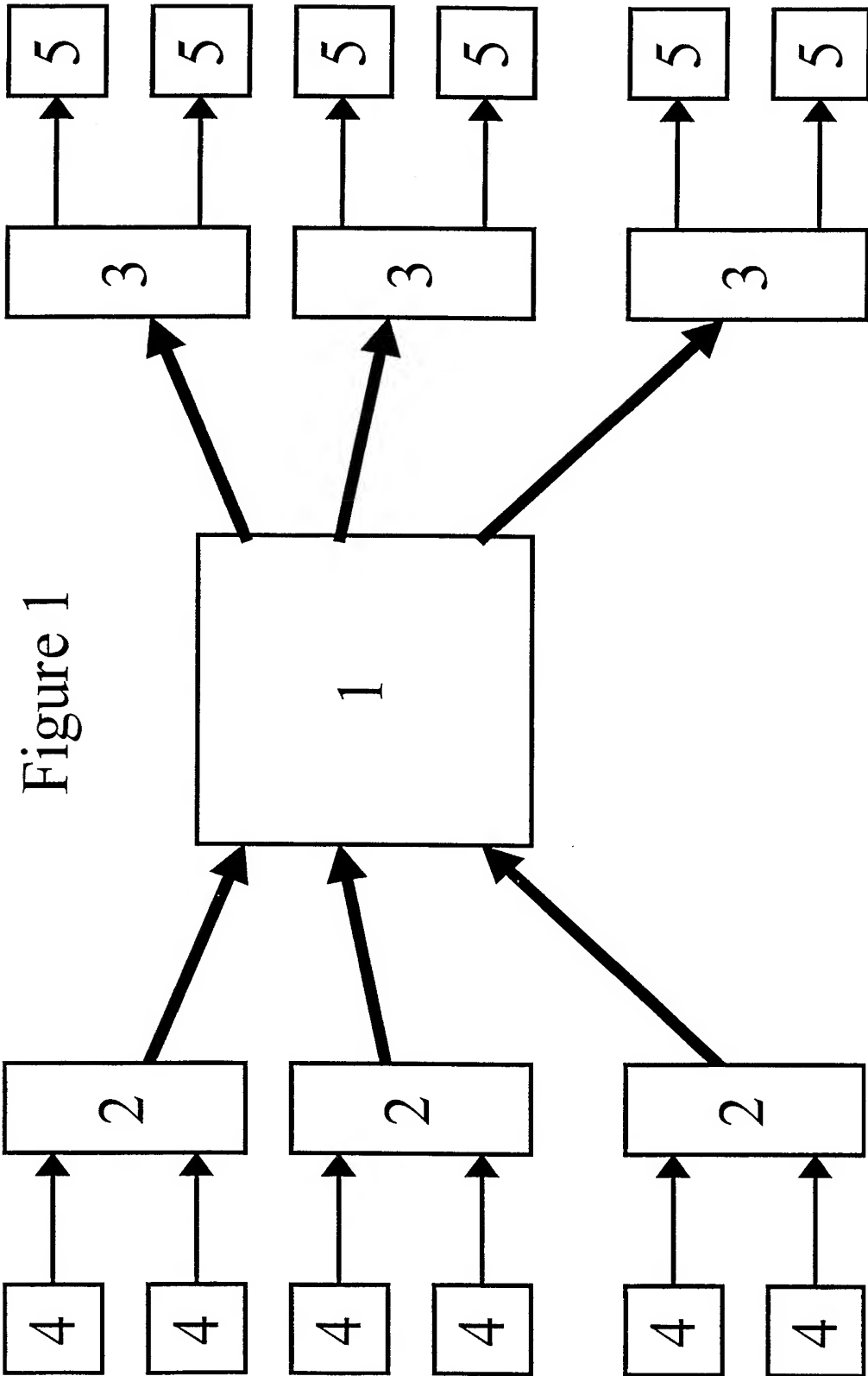


Figure 1

Figure 2

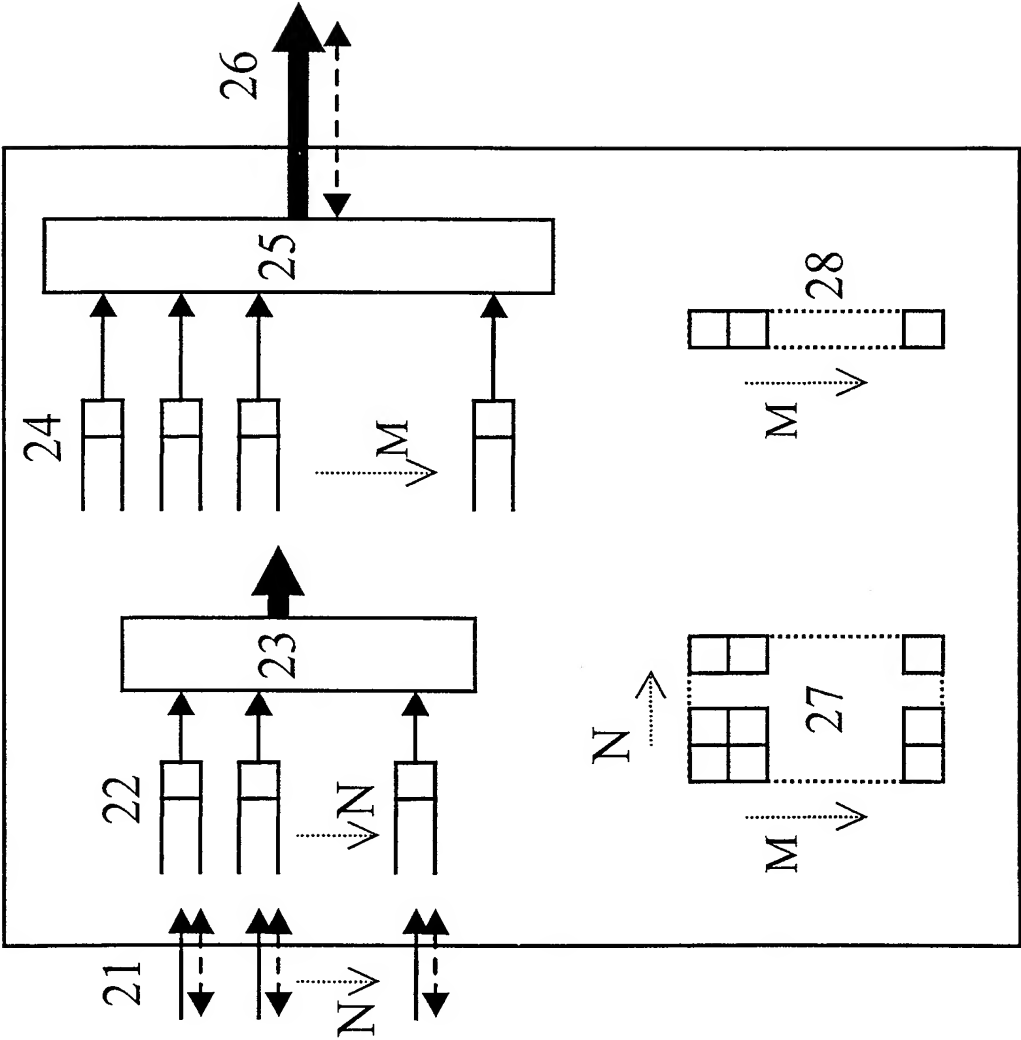


Figure 3

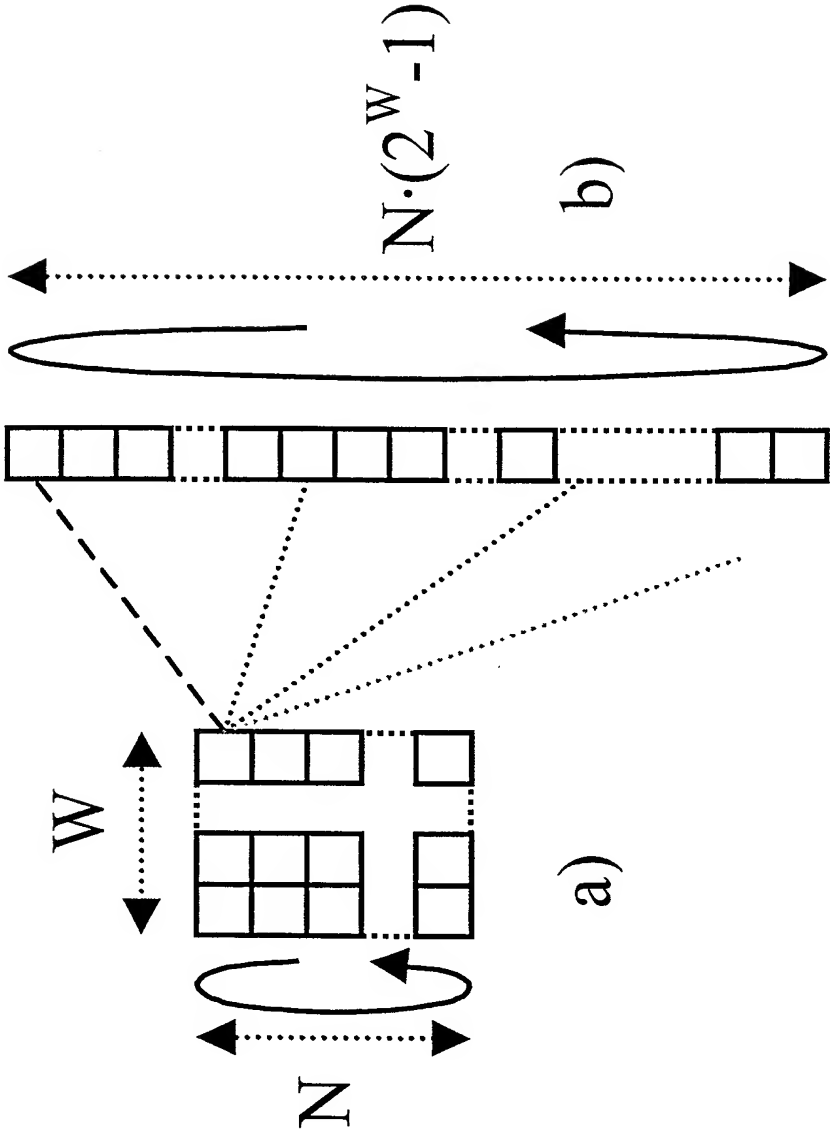


Figure 4

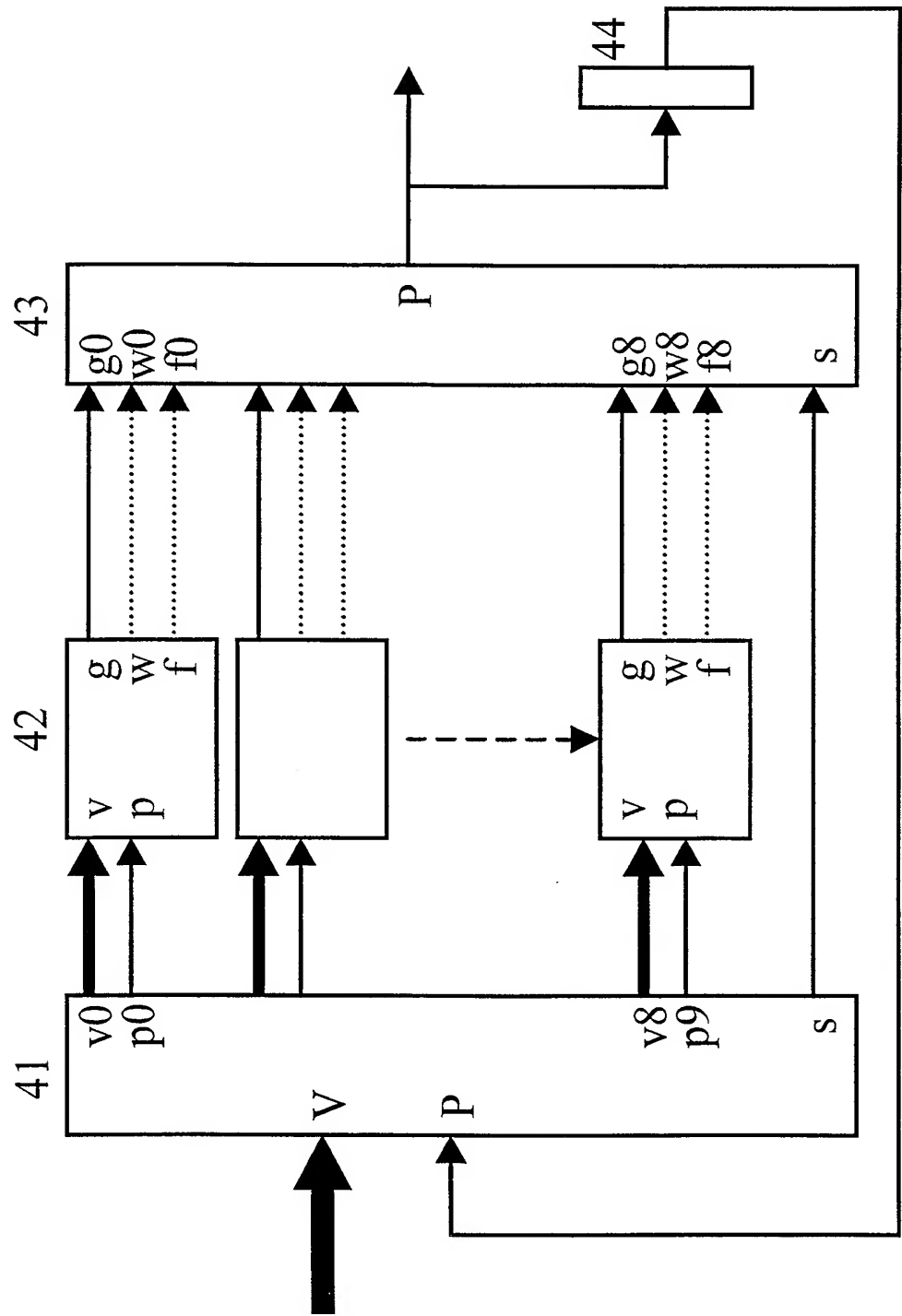


Figure 5

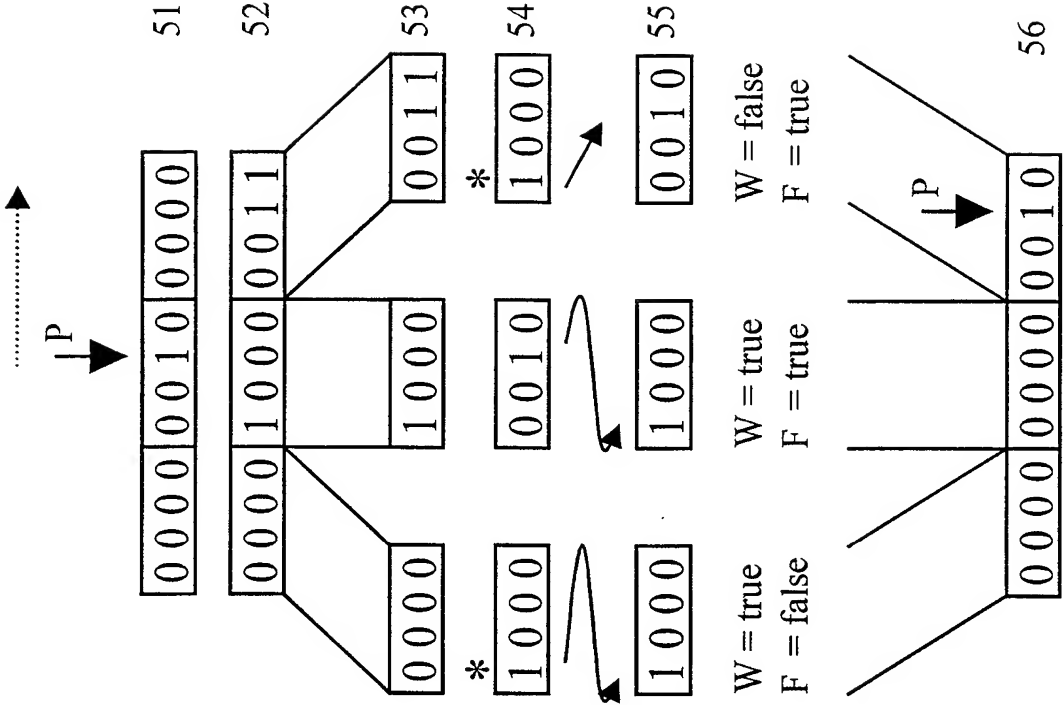
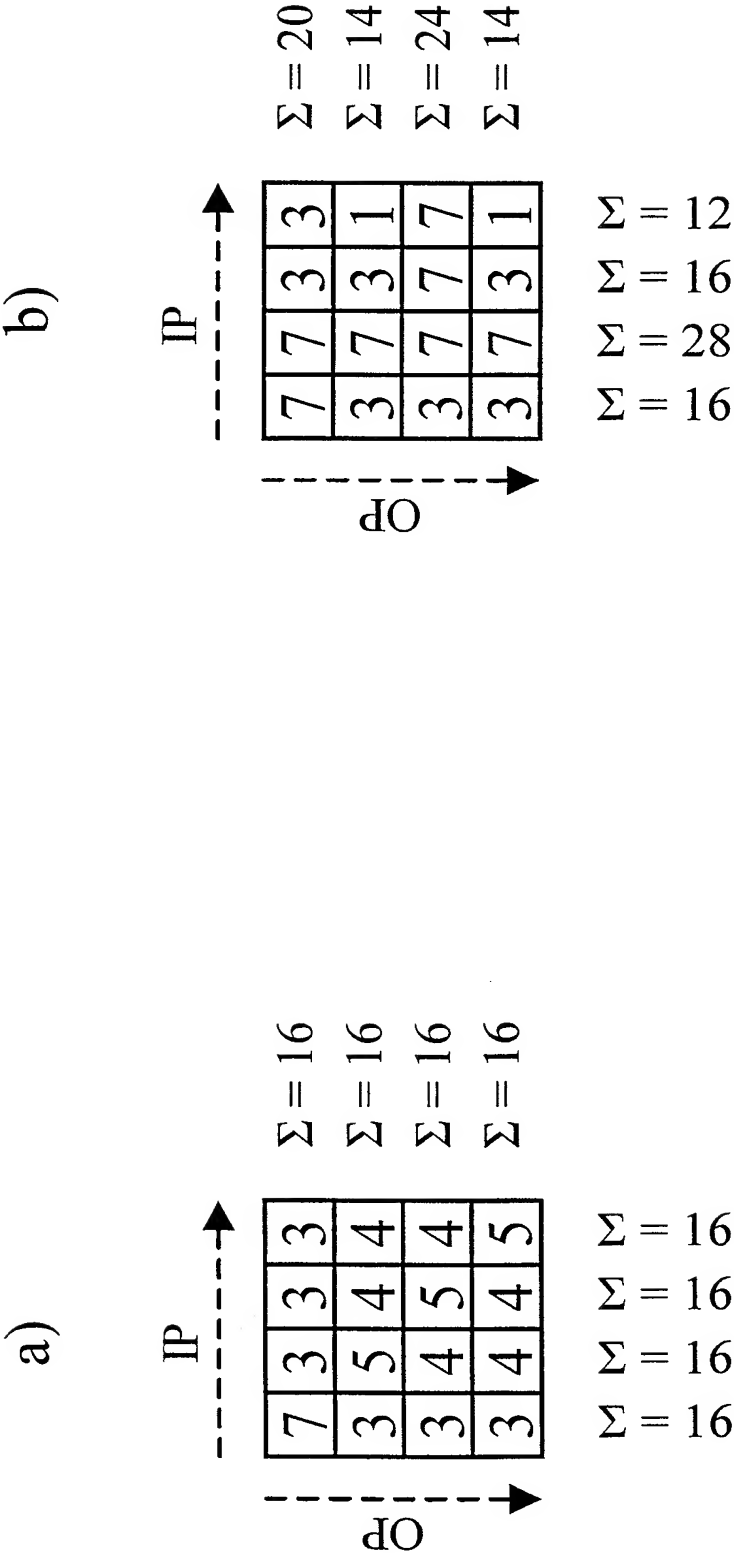
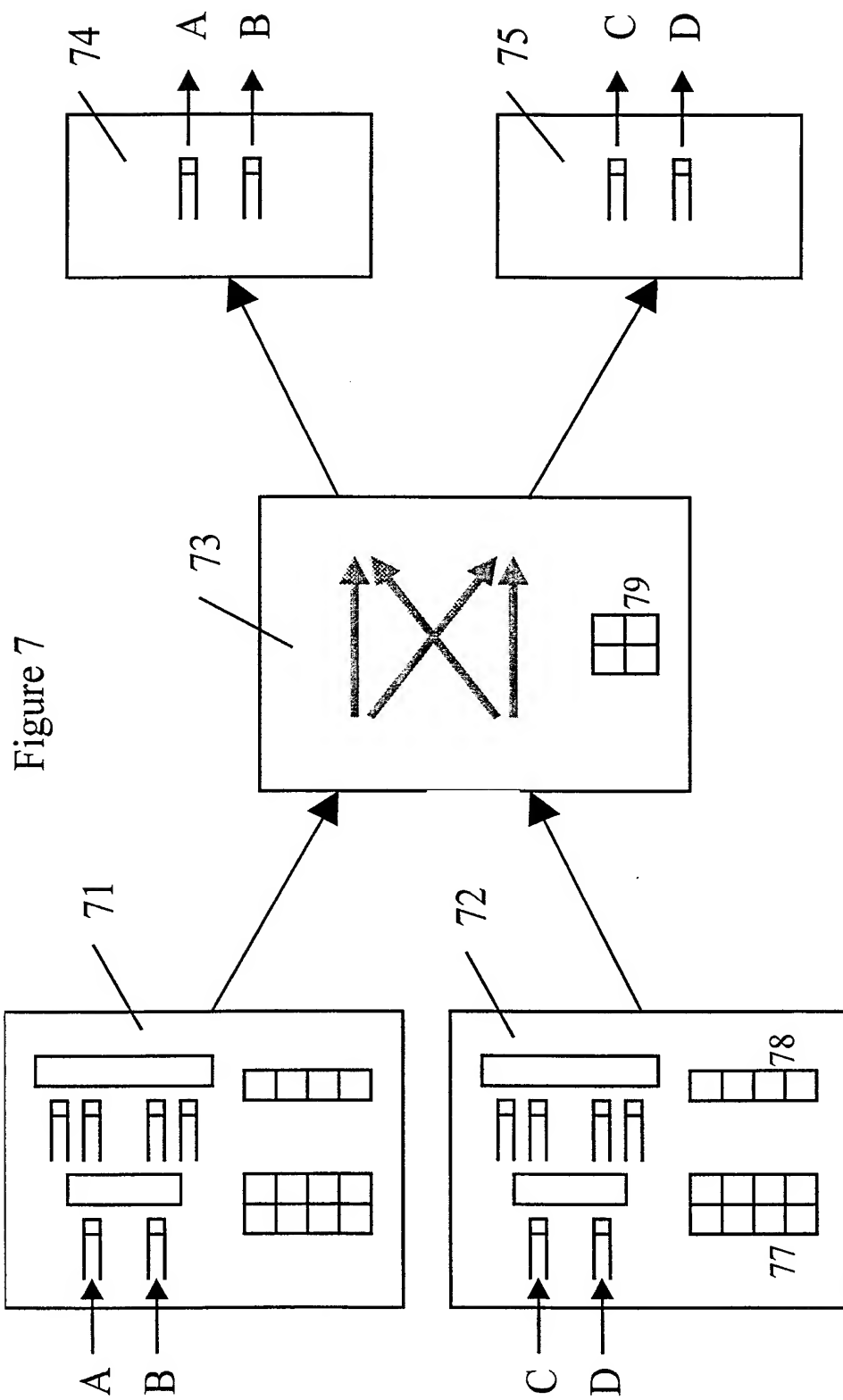


Figure 6





INTERNATIONAL SEARCH REPORT

International Application No
PCT/GB 99/04007

A. CLASSIFICATION OF SUBJECT MATTER
IPC 7 H04L12/56

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
IPC 7 H04L

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|------------|--|-----------------------|
| X | WALLMEIER E ET AL: "TRAFFIC CONTROL IN ATM SWITCHES WITH LARGE BUFFERS" ITC SPECIALISTS SEMINAR, NL, LEIDSCHENDAM, KPN RESEARCH, vol. SEMINAR 9, 1995, pages 45-60, XP000683145 abstract Figure 3-3 page 48, line 18-23 --- | 1-8 |
| A | WO 96 21303 A (STRATACOM INC) 11 July 1996 (1996-07-11) page 3-4 --- -/-- | 3,4 |

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

Date of the actual completion of the international search

17 March 2000

Date of mailing of the international search report

28/03/2000

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Dhondt, E

INTERNATIONAL SEARCH REPORT

International Application No

PCT/GB 99/04007

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|------------|---|-----------------------|
| A | <p>HUI ZHANG ET AL: "COMPARISON OF RATE-BASED SERVICE DISCIPLINES" COMPUTER COMMUNICATIONS REVIEW,US,ASSOCIATION FOR COMPUTING MACHINERY. NEW YORK, vol. 21, no. 4, 1 September 1991 (1991-09-01), pages 113-121, XP000234931 ISSN: 0146-4833 page 115, left-hand column, line 3 -right-hand column, line 5 -----</p> | 1-8 |

INTERNATIONAL SEARCH REPORT

Information on patent family members

International: plication No

PCT/GB 99/04007

| Patent document cited in search report | Publication date | Patent family member(s) | Publication date |
|---|---------------------|----------------------------|---------------------|
| WO 9621303 A | 11-07-1996 | US 5561663 A | 01-10-1996 |
| | | AU 4409596 A | 24-07-1996 |
| <hr/> | | | |